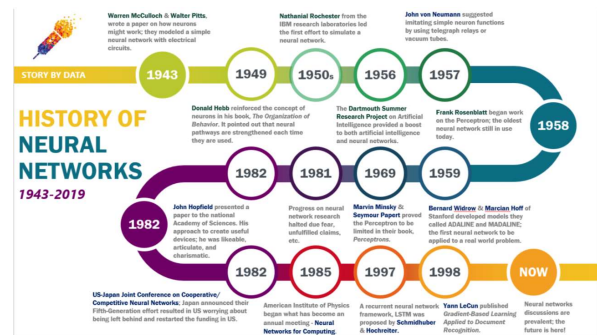


Unit – III: NEURAL NETWORKS

- ❖ Characteristics of Neural Networks, Historical Development of Neural Networks Principles. Artificial Neural Networks: Terminology, Models of Neuron, Topology, Basic Learning Laws, Pattern Recognition Problem, Basic Functional Units, Pattern Recognition Tasks by the Functional Units.

154

Historical Development of Neural Networks



155

Artificial Intelligence, Machine Learning and Deep Learning

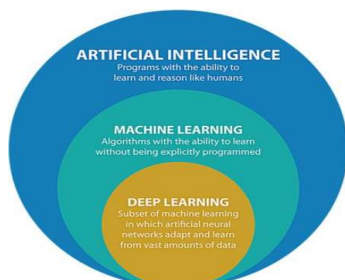
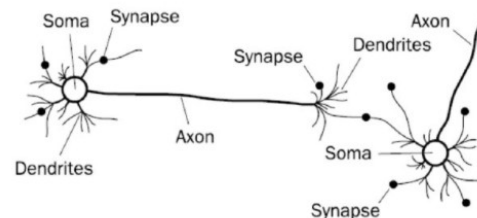


Image Source: <https://datacatchup.com/artificial-intelligence-machine-learning-and-deep-learning/>

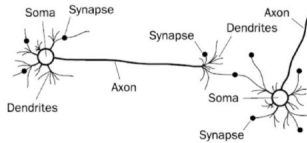
156

Biological Neural Networks

- A nerve cell neuron is a special biological cell that processes information.



157



Dendrites – They are tree-like branches, responsible for receiving the information from other neurons it is connected to. In other sense, we can say that they are like the ears of neuron.

Soma – It is the cell body of the neuron and is responsible for processing of information, they have received from dendrites.

Axon – It is just like a cable through which neurons send the information.

Synapses – It is the connection between the axon and other neuron dendrites.

158

Characteristic of Biological Neuron

- 1) Non-linearity
- 2) Input-output Mapping
- 3) Adaptability
- 4) Evidential Response (degree of 'confidence')
- 5) Fault Tolerance
- 6) Massively Parallel Computing

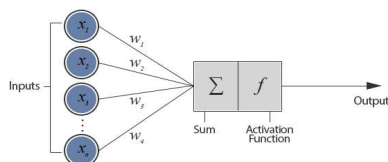
159

Artificial Neural Networks (ANN)

➤ Artificial neural networks, usually simply called neural networks, are **computing systems (mathematical function)**, inspired by the **biological neural networks**.

➤ **Artificial neurons** are elementary units in an **artificial neural network**.

➤ The weights can either amplify or deamplify the original input signal. For example, if the input is 1 and the input's weight is 0.2 the input will be decreased to 0.2.



160

Most artificial neurons have three things –

- an input transformation function
- an activation function
- an output transformation function

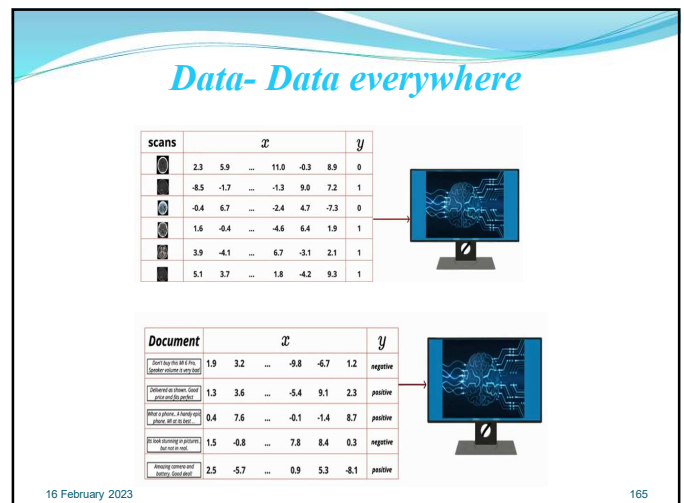
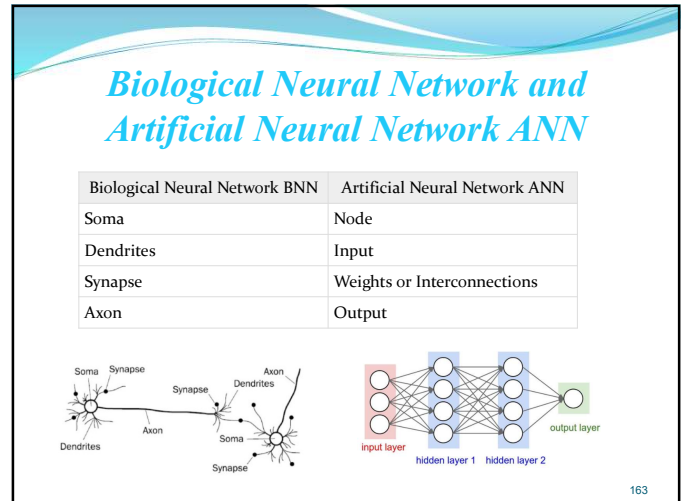
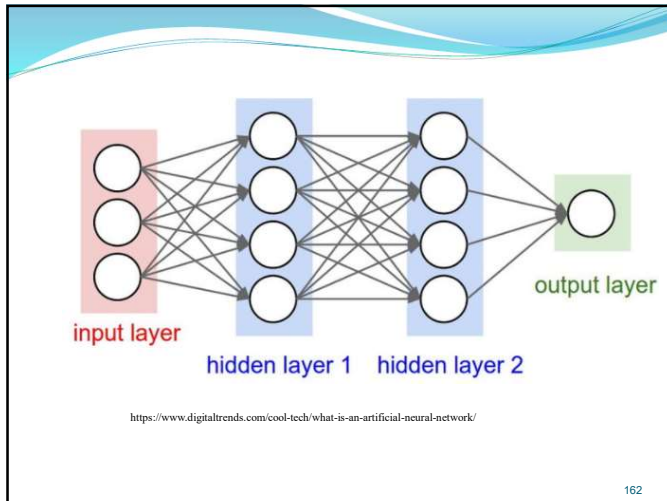
❑ Input transformation functions process the incoming signal – usually multiply and sum Artificial Neurons


❑ Activation functions process the input signal

- Threshold Function
- Linear
- Saturated linear
- Sigmoid

❑ Output functions process the outgoing signal – usually linear

161






1.3	-4.3	2.1	-6.7	...	1.5	8.9	10.1	-4.5
2.6	7.9	-0.3	8.1	...	-4.2	0.3	1.2	9.4
-5.2	-3.2	4.2	0.3	...	3.5	8.3	-1.4	-8.7
2.3	-5.6	-1.2	7.8	...	9.9	10.1	-1.1	3.5
8.5	2.1	-6.3	5.3	...	7.2	-1.3	-4.5	11.8


- ✓ All data encoded as numbers
- ✓ Typically high dimensional \mathbb{R}^n

16 February 2023 166

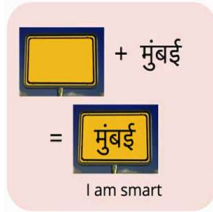
Data Curation



I am lucky



I am rich



I am smart

16 February 2023 167

Data Augmentation

- Data augmentation is an integral process in deep learning, as in deep learning we need large amounts of data and in some cases it is not feasible to collect thousands or millions of images, so data augmentation comes to the rescue.
- It helps us to increase the size of the dataset and introduce variability in the dataset.
- Data augmentation is the process of increasing the amount and diversity of data.
- We do not collect new data, rather we transform the already present data.

16 February 2023 168

Operations in Data Augmentation

- 1) Rotation:** Rotation operation as the name suggests, just rotates the image by a certain specified degree.
- 2) Zooming:** Zooming operation allows us to either zoom in or zoom out.
- 3) Cropping:** Cropping allows us to crop the image or select a particular area from an image.
- 4) Flipping:** Flipping allows us to flip the orientation of the image. We can use horizontal or vertical flip.
- 5) Changing the brightness level:** This feature helps us to combat illumination changes.

You can encounter a scenario where most of your dataset comprises of images having a similar brightness level e.g. collecting the images of employees entering the office, by augmenting the images we make sure that our model is robust and is able to detect the person even in different surroundings.

16 February 2023 169

What you want to do with data? Task

a Hello John,

Input

Output

Task (From description to structure specification)

a Hello John,

Input

Output

16 February 2023

171

Task (From description + review to write FAQs)

a Hello John,

Input

Output

16 February 2023

172

Task (From description + review + FAQs to Question Answering)

a Hello John,

Input

Output

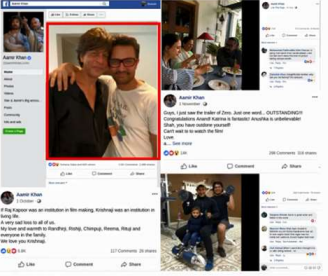
Can I shoot high quality videos?
Yes, you can shoot 4K HD videos on the device.

16 February 2023


173

Task (From image identify people)

Input



Output




16 February 2023

174

Task (From image identify activity)

Input



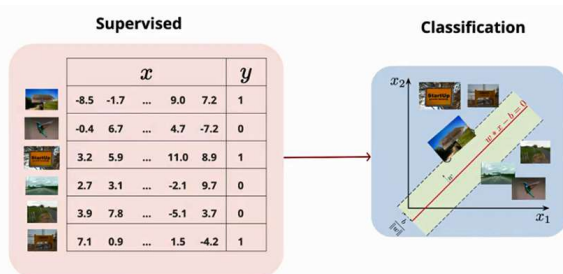
Output



16 February 2023

175

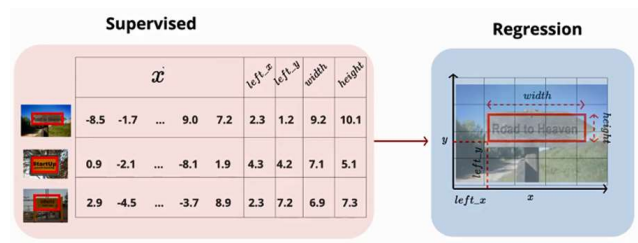
Task (Classification)



16 February 2023

176

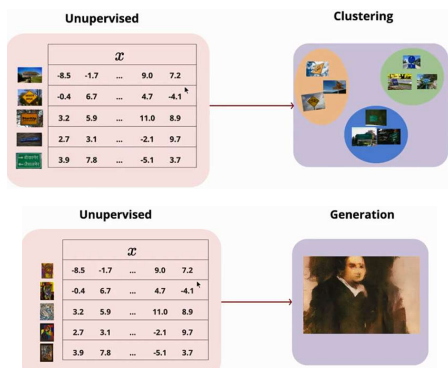
Task (Regression)



16 February 2023

177

Task (Clustering)

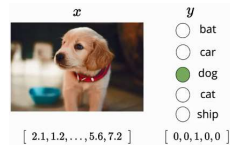


16 February 2023

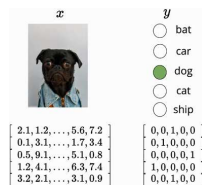
178

What is the mathematical formulation of a task?

Models



$$y = f(x) \text{ [true relation, unknown]}$$



$$y = f(x) \text{ [true relation, unknown]}$$

$$\hat{y} = \hat{f}(x) \text{ [our approximation]}$$

16 February 2023

179

What are the choice of \hat{f} ?

Models

$$\hat{y} = \hat{f}(x) \text{ [our approximation]}$$

x	y
0.5	14.8
0.2	13.3
0.6	11.6
...	...
0.3	6.16

$$\hat{y} = mx + c$$

Data is drawn from the following function

$$\hat{y} = \hat{f}(x) \text{ [our approximation]}$$

x	y
0.5	14.8
0.2	13.3
0.6	11.6
...	...
0.3	6.16

$$\hat{y} = mx + c$$

$$\hat{y} = ax^2 + bx + c$$

Data is drawn from the following function

16 February 2023

180

What are the choice of \hat{f} ?

Models

$$\hat{y} = \hat{f}(x) \text{ [our approximation]}$$

x	y
0.5	14.8
0.2	13.3
0.6	11.6
...	...
0.3	6.16

$$\hat{y} = mx + c$$

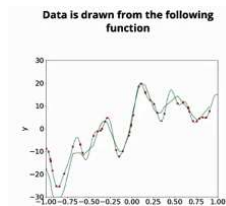
$$\hat{y} = ax^2 + bx + c$$

$$\hat{y} = ax^3 + bx^2 + cx + d$$

$$\hat{y} = ax^4 + bx^3 + cx + d$$

$$\vdots$$

$$\hat{y} = ax^{25} + bx^{24} + \dots + cx + d$$



16 February 2023

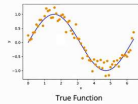
181

How do we know which model is better?

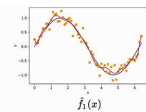
Loss function

- A loss function is a function that compares the target and predicted output values; measures how well the neural network models the training data. When training, we aim to minimize this loss between the predicted and target outputs.

	x	y
1	0.00	0.24
2	0.10	0.08
3	0.20	0.12
....
n	6.40	0.36



	x	y	$\hat{f}_1(x)$
1	0.00	0.24	0.25
2	0.10	0.08	0.09
3	0.20	0.12	0.11
....
n	6.40	0.36	0.36



16 February 2023

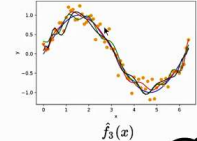
$$\hat{f}_1(x) = 1.79x^{25} - 4.54x^{24} + \dots - 1.48x + 2.48$$

182

How do we know which model is better?

Loss function

	x	y	$\hat{f}_1(x)$	$\hat{f}_2(x)$	$\hat{f}_3(x)$
1	0.00	0.24	0.25	0.32	0.08
2	0.10	0.08	0.09	0.30	0.20
3	0.20	0.12	0.11	0.31	0.14
....
n	6.40	0.36	0.36	0.22	0.15



$$\hat{f}_1(x) = 1.79x^{25} - 4.54x^{24} + \dots - 1.48x + 2.48$$

$$\hat{f}_2(x) = 2.27x^{25} + 9.89x^{24} + \dots + 2.79x + 3.22$$

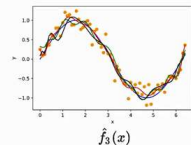
$$\hat{f}_3(x) = 3.78x^{25} + 1.57x^{24} + \dots + 1.01x + 8.68$$



16 February 2023

183

	x	y	$\hat{f}_1(x)$	$\hat{f}_2(x)$	$\hat{f}_3(x)$
1	0.00	0.24	0.25	0.32	0.08
2	0.10	0.08	0.09	0.30	0.20
3	0.20	0.12	0.11	0.31	0.14
....
n	6.40	0.36	0.36	0.22	0.15



$$\hat{f}_1(x) = 1.79x^{25} - 4.54x^{24} + \dots - 1.48x + 2.48$$

$$\hat{f}_2(x) = 2.27x^{25} + 9.89x^{24} + \dots + 2.79x + 3.22$$

$$\hat{f}_3(x) = 3.78x^{25} + 1.57x^{24} + \dots + 1.01x + 8.68$$

$$\mathcal{L}_1 = \sum_{i=1}^n (y_i - \hat{f}_1(x_i))^2$$

$$\mathcal{L}_2 = \sum_{i=1}^n (y_i - \hat{f}_2(x_i))^2$$

- ✓ Square Error Loss
- ✓ Cross Entropy Loss
- ✓ KL divergence

16 February 2023

184

	x	y	$\hat{f}_1(x)$	$\hat{f}_2(x)$	$\hat{f}_3(x)$
1	0.00	0.24	0.25	0.32	0.08
2	0.10	0.08	0.09	0.30	0.20
3	0.20	0.12	0.11	0.31	0.14
....
n	6.40	0.36	0.36	0.22	0.15

$$\begin{aligned} \mathcal{L}_1 &= \sum_{i=1}^n (y_i - \hat{f}_1(x_i))^2 \\ &= (0.24 - 0.25)^2 + (0.08 - 0.09)^2 + (0.12 - 0.11)^2 + \dots + (0.36 - 0.36)^2 \\ &= 1.38 \end{aligned}$$

$$\mathcal{L}_1 = \sum_{i=1}^n (y_i - \hat{f}_1(x_i))^2 = 1.38$$

$$\mathcal{L}_2 = \sum_{i=1}^n (y_i - \hat{f}_2(x_i))^2 = 2.02$$

$$\mathcal{L}_3 = \sum_{i=1}^n (y_i - \hat{f}_3(x_i))^2 = 2.34$$

16 February 2023

185

How do we identify the parameter of the model?

Learning Algorithm (Optimizer)

- The weights and bias are modified using a function called Optimization Function.
- Optimization functions usually calculate the gradient i.e. the **partial derivative of loss function with respect to weights**, and the weights are modified in the opposite direction of the calculated gradient. This cycle is repeated until we reach the minima of loss function.

Budget in (100crs)	Box Office Collection in (100 crs)	Action Scene in times (100 mins)	IMDB Rating
0.55	0.66	0.22	4.8
0.68	0.91	0.77	7.2
0.66	0.88	0.67	6.7
...
0.72	0.94	0.97	8.1
0.58	0.74	0.35	5.3

$$\hat{f}(x) = ax_1^2 + bx_2^3 + cx_3^2$$

$$\hat{f}(x) = 3.5x_1^2 + 2.5x_2^3 + 1.2x_3^2$$

$$\mathcal{L}_1 = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

16 February 2023

86

$$\hat{f}(x) = ax_1^2 + bx_2^3 + cx_3^2$$

Budget (100crs)	Box Office Collection(100 crs)	Action Scene times (100 mins)	IMDB Rating
0.55	0.66	0.22	4.8
0.68	0.91	0.77	7.2
0.66	0.88	0.67	6.7
...
0.72	0.94	0.97	8.1
0.58	0.74	0.35	5.3

$$\mathcal{L}_1 = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

✗ In practice, brute force search is infeasible

Budget (100crs)	Box Office Collection(100 crs)	Action Scene times (100 mins)	IMDB Rating
0.55	0.66	0.22	4.8
0.68	0.91	0.77	7.2
0.66	0.88	0.67	6.7
...
0.72	0.94	0.97	8.1
0.58	0.74	0.35	5.3

$$\hat{f}(x) = ax_1^2 + bx_2^3 + cx_3^2$$

Gradient Descent ++
Adagrad
RMSProp
Adam

Many optimization solvers are available






$$\min_{a,b,c} \mathcal{L}_1 = \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

16 February 2023

187

How do we score for our model?

Evaluation






	Top - 1	
	True Labels	Predicted Labels
	[2.1, 1.2, ..., 5.6, 7.8]	1 ✓
	[3.5, 6.6, ..., 2.5, 6.3]	2 3 ✗
	[6.3, 2.6, ..., 4.5, 3.8]	3 1 ✗
	[2.8, 3.6, ..., 7.5, 2.1]	4 4 ✓
	[2.2, 1.7, ..., 2.5, 1.8]	5 5 ✓
	[6.3, 2.6, ..., 4.5, 3.8]	3 2 ✗
	[1.9, 3.3, ..., 4.2, 1.1]	5 5 ✓

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{4}{7} = 0.55$$

16 February 2023

188

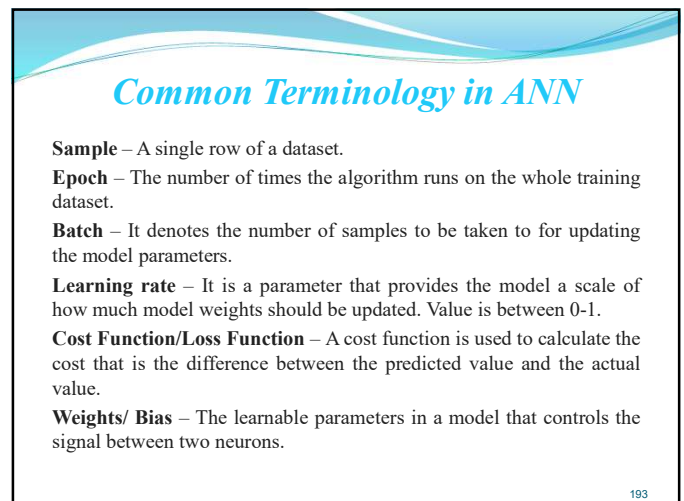
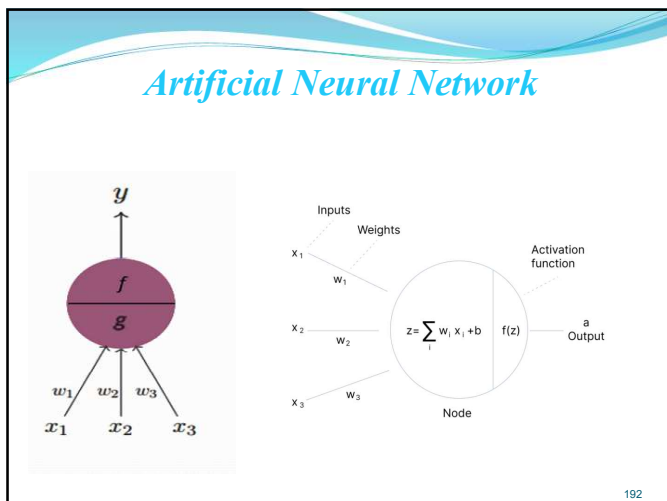
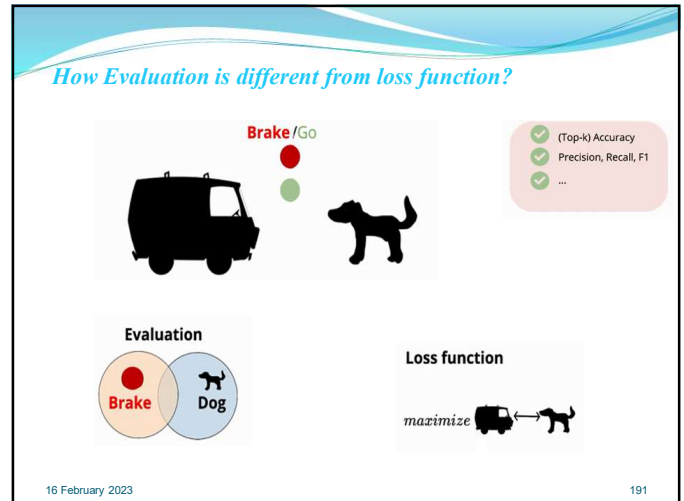
Top - 3

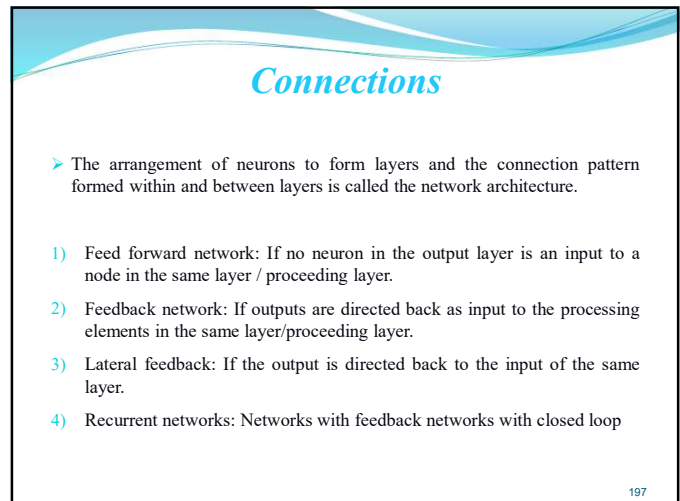
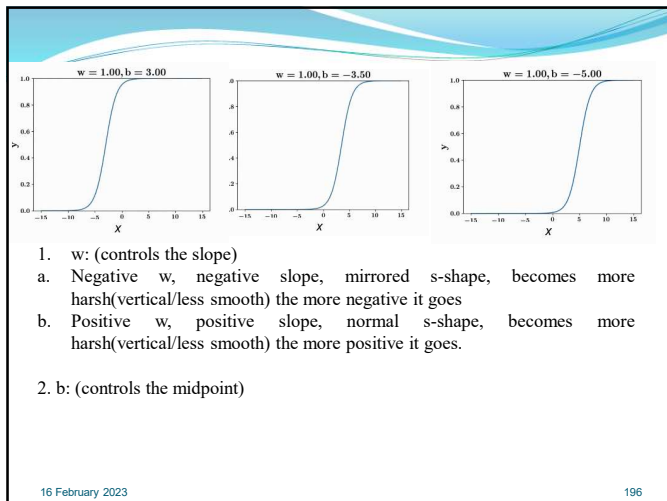
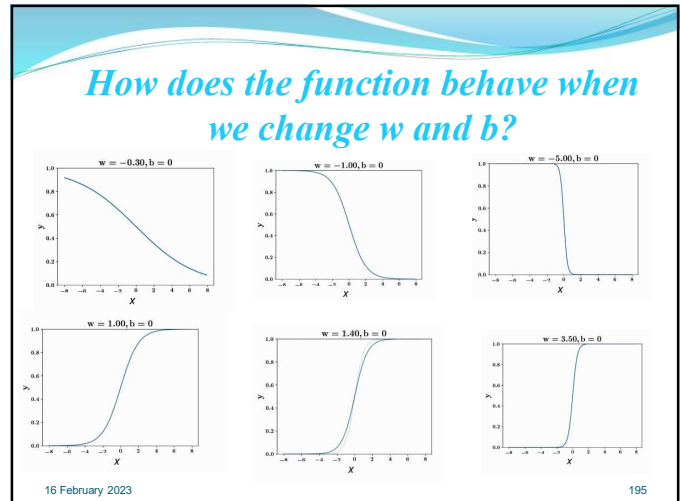
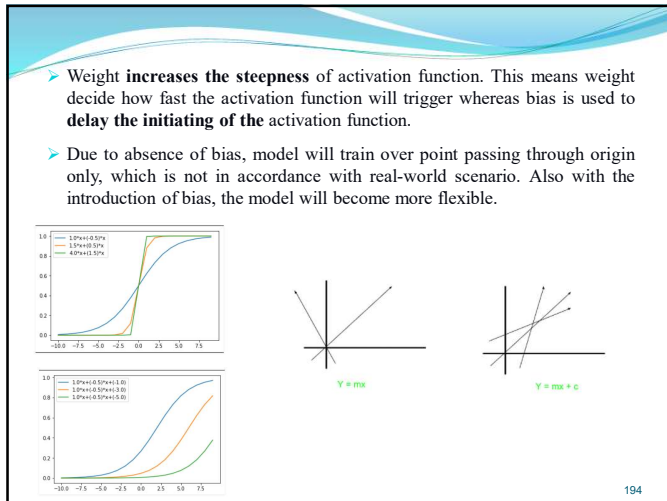
	True Labels	Predicted Labels	
	[2.1, 1.2, ..., 5.6, 7.8]	1 [1, 2, 3]	✓
	[3.5, 6.6, ..., 2.5, 6.3]	2 [1, 2, 3]	✓
	[6.3, 2.6, ..., 4.5, 3.8]	3 [1, 2, 3]	✓
	[2.8, 3.6, ..., 7.5, 2.1]	4 [4, 5, 3]	✓
	[2.2, 1.7, ..., 2.5, 1.8]	5 [5, 2, 1]	✓
	[6.3, 2.6, ..., 4.5, 3.8]	3 [2, 1, 4]	✗
	[1.9, 3.3, ..., 4.2, 1.1]	5 [5, 4, 1]	✓

$$\text{Accuracy} = \frac{\text{Number of correct predictions in top-3}}{\text{Total number of predictions}} = \frac{6}{7} = 0.86$$

16 February 2023

189

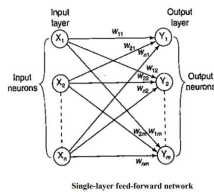




➤ There exist five basic types of connection architecture.

- 1) Single layer feed forward network
- 2) Multilayer feed-forward network
- 3) Single node with its own feedback
- 4) Single-layer recurrent network
- 5) Multilayer recurrent network

Single layer feed forward network: In this type of network, we have only two layers input layer and output layer but the input layer does not count because no computation is performed in this layer. The output layer is formed when different weights are applied on input nodes and the cumulative effect per node is taken. After this, the neurons collectively give the output layer to compute the output signals.



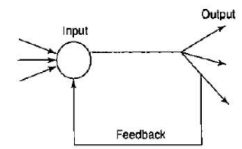
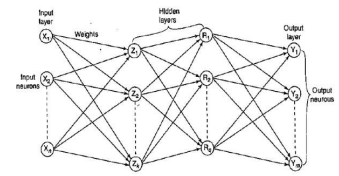
198

Multilayer feed-forward network:

This network is formed by the interconnection of several layers. Input layer receives input and buffers input signal. Output layer generated output.

Layer between input and output is called hidden layer. Hidden layer is internal to the network. There are Zero to several hidden layers in a network.

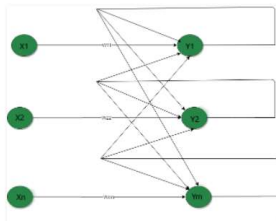
More the hidden layer more is the complexity of network, but efficient output is produced.



Single node with own feedback

Single-layer recurrent network: This network is a single-layer network with a feedback connection in which the processing element's output can be directed back to itself or to another processing element or both.

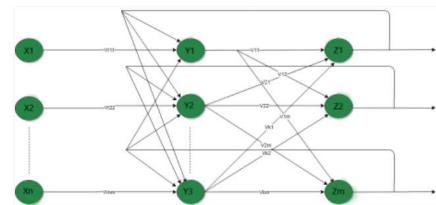
This allows it to exhibit dynamic temporal behavior for a **time sequence**. Unlike feedforward neural networks, RNNs can use their internal state (memory) to process sequences of inputs.



200

Multi-layer recurrent network: In this network, processing element output can be directed to the processing element in the same layer and in the preceding layer forming a multilayer recurrent network.

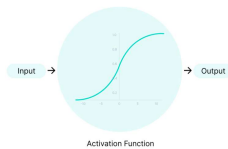
They perform the same task for every element of a sequence, with the output being dependent on the previous computations. Inputs are not needed at each time step. The main feature of a Recurrent Neural Network is its hidden state, which captures some information about a sequence.



201

Activation Function

- It will help to determine whether the neuron will **fire or not**.
- An activation function is to **add non-linearity** to the neural network.
- If we have a neural network working **without the activation** functions. In that case, **every neuron will only be performing a linear transformation** on the inputs using the weights and biases. It's because it **doesn't matter how many hidden layers we attach in the neural network**; all layers will behave in the same way because the composition of two linear functions is a linear function itself.



202

Activation Function Types

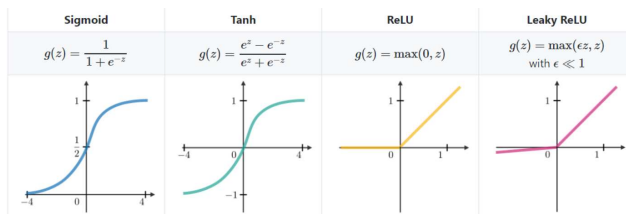
- 1) Linear Function
- 2) Binary Step Function
- 3) Non-Linear Function

Linear Function

- The linear activation function is also known as *Identity Function* where the **activation is proportional to the input**.

203

Non-linear Activation function



204

Softmax Function

- The softmax function is a function that turns a **vector of K real values into a vector of K real values that sum to 1**.
- The input values can be positive, negative, zero, or greater than one, but the softmax transforms them into values between 0 and 1, so that they can be interpreted as probabilities.
- If one of the inputs is small or negative, the softmax turns it into a small probability, and if an input is large, then it turns it into a large probability, but it will always remain between 0 and 1.

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

205

$$\begin{bmatrix} 8 \\ 5 \\ 0 \end{bmatrix}$$

$$e^{z_1} = e^8 = 2981.0$$

$$e^{z_2} = e^5 = 148.4$$

$$e^{z_3} = e^0 = 1.0$$

$$\sum_{j=1}^K e^{z_j} = e^{z_1} + e^{z_2} + e^{z_3} = 2981.0 + 148.4 + 1.0 = 3130.4$$

$$\sigma(\vec{z})_1 = \frac{2981.0}{3130.4} = 0.9523$$

$$\sigma(\vec{z})_2 = \frac{148.4}{3130.4} = 0.0474$$

$$\sigma(\vec{z})_3 = \frac{1.0}{3130.4} = 0.0003$$

206

Choosing the Right One

- 1) Use the ReLu or LeakyRelu function in hidden layers only
- 2) In binary classification remember to use the Sigmoid function in the output layer
- 3) In multi-class classification (when classes to predict are more than 2) problem use Softmax function in the output layer
- 4) Due to the vanishing gradient problem 'Sigmoid' and 'Tanh' activation functions are avoided sometimes in deep neural network architectures
- 5) Always remember you can also invent your **own activation** function and can check its performance.

207

Learning Rules

- Learning rule or Learning process is a method or a mathematical logic. It improves the Artificial Neural Network's performance and applies this rule over the network.
- Thus learning rules updates the weights and bias levels of a network when a network simulates in a specific data environment.

- 1) Hebbian learning
- 2) Perceptron learning rule
- 3) Delta learning
- 4) Correlation learning rule
- 5) Outstar learning rule

208

Hebbian Learning Rule

- The Hebbian rule was the first learning rule. In 1949 Donald Hebb developed it as learning algorithm of the neural network.
- The Hebb learning rule assumes that – If two neighbor neurons **activated and deactivated** at the same time. Then the weight connecting these neurons should **increase**.
- For neurons operating in the **opposite phase**, the weight between them should **decrease**.
- If there is no signal correlation, the weight **should not change**.

$$w_i(\text{new}) = w_i(\text{old}) + x_i y$$

209

- When inputs of both the nodes are either **positive or negative**, then a **strong positive weight exists** between the nodes.
- If the input of a node is **positive and negative** for other, a strong **negative weight exists** between the nodes.
- At the start, values of all **weights are set to zero**. This learning rule can be used for both soft- and hard-activation functions.

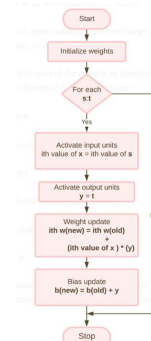
➤ Hebb's network is suited more for **bipolar data**. If binary data is used, the weight update formula cannot distinguish two conditions namely:

1. A training pair in which an input unit is "on" and the target value is "off".
2. A training pair in which both the input unit and the target value is "off".

210

Flow Diagram

s: t refers to each training input and target output pair. Till there exists a pair of training input and target output, the training process takes place; else it is stopped.



211

STEP 1: Initialize the weights and bias to '0' i.e $w_1=0, w_2=0, \dots, w_n=0$.

STEP 2: 3–5 have to be performed for each input training vector and target output pair i.e. s:t (s=training input vector, t=training output vector)

STEP 3: Input units activation are set and in most of the cases is an identity function (one of the types of an activation function) for the input layer;

$$x_i = s_i \text{ for } i=1 \text{ to } n$$

Identity Function: Its a linear function and defined as $f(x)=x$ for all x

STEP 4: Output units activations are set y:t

STEP 5: Weight adjustments and bias adjustments are performed;

$$w_i(\text{new}) = w_i(\text{old}) + (x_i * y)$$

$$\text{bias}(\text{new}) = \text{bias}(\text{old}) + y$$

212

Implementation of AND Gate with Hebbian Learning Rule

Set all weights to zero, $w_i = 0$ for $i=1$ to n , and **bias to zero**.

	INPUT			TARGET	
	x_1	x_2	b	y_1	y
X_1	-1	-1	1	Y_1	-1
X_2	-1	1	1	Y_2	-1
X_3	1	-1	1	Y_3	-1
X_4	1	1	1	Y_4	1

$$w_i(\text{new}) = w_i(\text{old}) + x_i y$$

$$b(\text{new}) = b(\text{old}) + y$$

Activation function used here is Bipolar Sigmoidal Function so the range is $[-1, 1]$.

213

Step 1:
Set weight and bias to zero, $w = [0\ 0\ 0]^T$ and $b = 0$.

Step 2:
Set input vector $X_i = S_i$ for $i = 1$ to 4.
 $X_1 = [-1\ -1\ 1]^T$
 $X_2 = [-1\ 1\ 1]^T$
 $X_3 = [1\ -1\ 1]^T$
 $X_4 = [1\ 1\ 1]^T$

Step 3:
Output value is set to $y = t$.

Step 4:
Modifying weights using Hebbian Rule:

$$\Delta w_1 = x_1 y = 1 \times 1 = 1$$

$$\Delta w_2 = x_2 y = 1 \times 1 = 1$$

$$\Delta b = y = 1$$

$$w_1(\text{new}) = w_1(\text{old}) + x_1 y$$

$$w_1(\text{new}) = w_1(\text{old}) + x_1 y = 0 + 1 \times 1 = 1$$

$$w_2(\text{new}) = w_2(\text{old}) + x_2 y = 0 + 1 \times 1 = 1$$

$$b(\text{new}) = b(\text{old}) + y = 0 + 1 = 1$$

The weight change here is

$$\Delta w_1 = x_1 y = 1 \times -1 = -1$$

$$\Delta w_2 = x_2 y = -1 \times -1 = 1$$

$$\Delta b = y = -1$$

The new weights here are

$$w_1(\text{new}) = w_1(\text{old}) + \Delta w_1 = 1 - 1 = 0$$

$$w_2(\text{new}) = w_2(\text{old}) + \Delta w_2 = 1 + 1 = 2$$

$$b(\text{new}) = b(\text{old}) + \Delta b = 1 - 1 = 0$$

Inputs				Weight Changes			Weights		
x1	x2	b	y	Δw_1	Δw_2	Δb	w1	w2	b
1	1	1	1				(0	0	0)
1	-1	1	-1						
-1	1	1	-1						
-1	-1	1	-1						

214

The separating line equation is given by

$$x_2 = \frac{-w_1}{w_2} x_1 - \frac{b}{w_2}$$

For all inputs, use the final weights obtained for each input to obtain the separating line. For the first input $[1\ 1\ 1]$, the separating line is given by

$$x_2 = \frac{-1}{1} x_1 - \frac{1}{1} \Rightarrow x_2 = -x_1 - 1$$

Similarly, for the second input $[1\ -1\ 1]$, the separating line is

$$x_2 = \frac{-0}{2} x_1 - \frac{0}{2} \Rightarrow x_2 = 0$$

For the third input $[-1\ 1\ 1]$, it is

$$x_2 = \frac{-1}{1} x_1 + \frac{1}{1} \Rightarrow x_2 = -x_1 + 1$$

Finally, for the fourth input $[-1\ -1\ 1]$, the separating line is

$$x_2 = \frac{-2}{2} x_1 + \frac{2}{2} \Rightarrow x_2 = -x_1 + 1$$

215

(A) First input

(B) Second input

(C) Third and fourth inputs

216

Perceptron Learning Rule

- This rule is an **error correcting** and a **supervised learning algorithm** of single **layer feedforward networks with linear activation function**.
- As being supervised in nature, to calculate the error, there would be a comparison between the desired/target output and the actual output.
- If there is any difference found, then a change must be made to the weights of connection.
- The perceptron network consists of three units, namely, sensory unit (input unit), associator unit (hidden unit), response unit (output unit).
- The input units are connected to hidden units with fixed weights having values **1, 0 or -1**, which are assigned at random.
- **The binary activation function** is used in hidden layer.
- The response unit has an activation of **1, 0 or -1**

217

The output of the perceptron network is given by

$$y = f(y_{in})$$

where $f(y_{in})$ is activation function and is defined as

$$f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} > \theta \\ 0 & \text{if } -\theta \leq y_{in} \leq \theta \\ -1 & \text{if } y_{in} < -\theta \end{cases}$$

If $y \neq t$, then

$$w_i(\text{new}) = w_i(\text{old}) + \alpha t x_i$$

$$b(\text{new}) = b(\text{old}) + \alpha t$$

"t" is +1 or -1 and α is the learning rate.

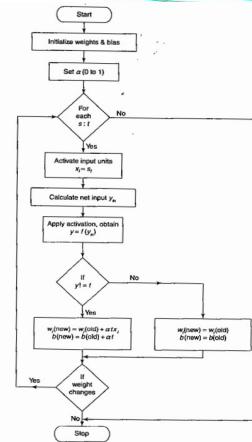
else,

$$w_i(\text{new}) = w_i(\text{old})$$

$$b(\text{new}) = b(\text{old})$$

- If no error occurs, there is no weight updation and hence the training process may be stopped.

218



219

Step 0: Initialize the weights and the bias (for easy calculation they can be set to zero). Also initialize the learning rate α ($0 < \alpha \leq 1$). For simplicity α is set to 1.

Step 1: Perform Steps 2-6 until the final stopping condition is false.

Step 2: Perform Steps 3-5 for each training pair indicated by x_i .

Step 3: The input layer containing input units is applied with identity activation functions

$$x_i = t_i$$

Step 4: Calculate the output of the network. To do so, first obtain the net input:

$$y_{in} = \theta + \sum_{i=1}^n x_i w_i$$

where "n" is the number of input neurons in the input layer. Then apply activation function over the net input calculated to obtain the output:

$$y = f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} > \theta \\ 0 & \text{if } -\theta \leq y_{in} \leq \theta \\ -1 & \text{if } y_{in} < -\theta \end{cases}$$

Step 5: Weight and bias adjustment: Compare the value of the actual (calculated) output and desired (target) output.

If $y \neq t$, then

$$w_i(\text{new}) = w_i(\text{old}) + \alpha x_i$$

$$b(\text{new}) = b(\text{old}) + \alpha t$$

else, we have

$$w_i(\text{new}) = w_i(\text{old})$$

$$b(\text{new}) = b(\text{old})$$

Step 6: Train the network until there is no weight change. This is the stopping condition for the network. If this condition is not met, then start again from Step 2.

220

Perceptron Summary

- It required initial weight to be assigned random values.
- Uses supervised learning
- Guarantees convergence in finite iterations
- Used in feedforward network

221

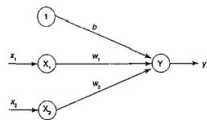
Implement AND function using perceptron networks for bipolar inputs and targets.

x_1	x_2	t
1	1	1
1	-1	-1
-1	1	-1
-1	-1	-1

$$y = f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} > 0 \\ 0 & \text{if } y_{in} = 0 \\ -1 & \text{if } y_{in} < 0 \end{cases}$$

$$\begin{aligned} \Delta w_1 &= \alpha x_1; \\ \Delta w_2 &= \alpha x_2; \\ \Delta b &= \alpha t \end{aligned}$$

$w_1 = w_2 = b = 0$ and $\theta = 0$. The learning rate α is set equal to 1.



$$\begin{aligned} w_1(\text{new}) &= w_1(\text{old}) + \alpha x_1 \\ w_1(\text{new}) &= w_1(\text{old}) + \alpha x_1 = 0 + 1 \times 1 \times 1 = 1 \\ w_2(\text{new}) &= w_2(\text{old}) + \alpha x_2 = 0 + 1 \times 1 \times 1 = 1 \\ b(\text{new}) &= b(\text{old}) + \alpha t = 0 + 1 \times 1 \times 1 = 1 \end{aligned}$$

• Calculate the net input

$$\begin{aligned} y_{in} &= b + x_1 w_1 + x_2 w_2 \\ &= 0 + 1 \times 0 + 1 \times 0 = 0 \end{aligned}$$

222

Input		Target	Net input	Calculated	Weight changes			Weights		
x_1	x_2	t	(y_{in})	(y)	Δw_1	Δw_2	Δb	w_1	w_2	b
EPOCH-1										
1	1	1	1					0	0	0
1	-1	1	-1							
-1	1	1	-1							
-1	-1	1	-1							

Input		Target	Net input	Calculated	Weight changes			Weights		
x_1	x_2	t	(y_{in})	(y)	Δw_1	Δw_2	Δb	w_1	w_2	b
EPOCH-1										
1	1	1	1	0	0	0	0	1	1	1
1	-1	1	-1	1	-1	1	-1	0	2	0
-1	1	1	-1	2	+1	-1	-1	1	1	-1
-1	-1	1	-1	-3	-1	0	0	0	1	-1
EPOCH-2										
1	1	1	1	1	0	0	0	1	1	-1
1	-1	1	-1	-1	-1	0	0	0	1	-1
-1	1	1	-1	-1	-1	0	0	0	1	-1
-1	-1	1	-1	-3	-1	0	0	0	1	-1

223

The final weights and bias after second epoch are

$$w_1 = 1, w_2 = 1, b = -1$$

Since the threshold for the problem is zero, the equation of the separating line is

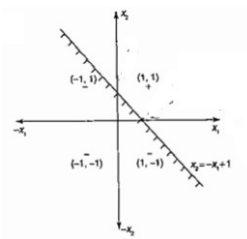
$$x_2 = -\frac{w_1}{w_2} x_1 - \frac{b}{w_2}$$

Here

$$\begin{aligned} w_1 x_1 + w_2 x_2 + b &> \theta \\ w_1 x_1 + w_2 x_2 + b &> 0 \end{aligned}$$

Thus, using the final weights we obtain

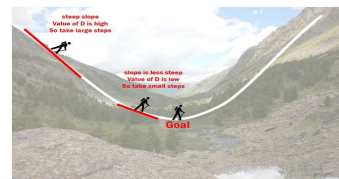
$$\begin{aligned} x_2 &= -\frac{1}{1} x_1 - \frac{(-1)}{1} \\ x_2 &= -x_1 + 1 \end{aligned}$$



224

Gradient Descent

- Gradient descent is an iterative optimization algorithm to find the minimum of a function. Here that function is our Loss Function.
- Loss function is the partial deviation with respect to w and b .



- He goes down the slope and takes large steps when the slope is steep and small steps when the slope is less steep.
- He decides his next position based on his current position and stops when he gets to the bottom of the valley which was his goal.

225

Delta Learning Rule (Widrow–Hoff Rule)

- The perceptron learning rule originates from the Hebbian assumption while the delta rule is derived from the gradient- descent method.
- The delta rule updates the weights between the connections so as to minimize the **difference between the net input to the output unit and the target value**.
- The major aim is to minimize all errors over all training patterns. This is done by reducing the error for each pattern, one at a time.
- Delta rule (DR) is similar to the Perceptron Learning Rule (PLR), with some differences:
 1. Error (Error (δ), in DR is not restricted to having values of 0, 1, or - 1 (as in PLR), but may have any value
 2. DR can be derived for any differentiable output/activation function f , whereas in PLR only works for threshold output function.

226

Case-I – when $t \neq y$, then

$$w(\text{new}) = w(\text{old}) + \Delta w$$

Case-II – when $t = y$, then

No change in weight

where

$$\Delta w_i = \eta(t - o)x_i$$

learning rate target value perceptron output input value

$$E(w_{ij}) = \frac{1}{2}(y_{\text{target}} - y_j)^2 \quad y_j = f(a_j) = \sum_i w_{ij} x_i$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial y_j} \cdot \frac{\partial y_j}{\partial w_{ij}} = -(y_{\text{target}} - y_j)x_i = -\delta x_i$$

$$w_i(\text{new}) = w_i(\text{old}) + \alpha(t - y_{ia})x_i$$

$$b(\text{new}) = b(\text{old}) + \alpha(t - y_{ia})$$

227